

第5卷第4期
2006年8月

江南大学学报(自然科学版)
Journal of Southern Yangtze University(Natural Science Edition)

Vol. 5 No. 4
Aug. 2006

文章编号:1671-7147(2006)04-0497-03

基于支持向量回归的钢材力学性能模型及应用

林豫柏, 罗键*
(厦门大学自动化系, 福建 厦门 361005)

摘要: 根据化学成分准确预测钢材产品的力学性能并及时调整相关生产的控制策略, 将有效地提高钢铁生产的最终产品质量. 支持向量机是一种建立在统计学习理论基础上的机器学习方法, 介绍了基于此算法基础上的一种 γ -支持向量回归机算法及其推导过程, 建立了基于 γ -支持向量回归机的钢材力学性能模型, 通过实际应用表明该模型比 Excel 回归预测具有更高的精度.

关键词: 预测; 支持向量机; 力学性能

中图分类号: TP 301

文献标识码: A

Model and Application Based on Supportable Recursive Vector for Forecasting Steel Mechanics Performance

LIN Yu-bai, LUO Jian*
(Department of Automation, Xiamen University, Xiamen 361005, China)

Abstract: The quality of final steel product will be raised efficiently by accurately forecasting the mechanics performances based on the chemical components and adjusting related control strategy. SVM is a kind of machine learning way based on statistics theory. A γ -SVM method based on this arithmetic and its deduce are introduced. Then, a model of steel mechanics performances based on γ -SVM is constructed. Finally, an application of this model through practical data indicates that γ -SVM get a higher precision than Excel recursive forecast.

Key words: forecast; SVM; mechanics performance

支持向量机(support vector machines, SVM)是 AT & T Bell 实验室的 Vapnik 于 20 世纪 90 年代提出来的, 是一种建立在统计学习理论基础上的机器学习方法^[1-2]. 它在解决小样本学习、非线性以及高维模式识别等问题中表现出许多特有的优势, 其基本思想可概括为: 首先通过非线性变换将输入空间变换到一个高维空间, 然后在这个新空间中求取最优线性分类

面, 而这种非线性变换是通过定义适当的内积函数实现的. 根据结构风险最小化准则, 在训练样本分类误差极小化的前提下, 尽量提高分类器的泛化推广能力. 从实施的角度看, 训练支持向量机等价于解一个线性约束的二次规划问题, 使得分隔特征空间中两类模式点的两个超平面之间距离最大, 而且它能保证得到的解为全局最优点. 训练支持向量机已广泛应用于

收稿日期: 2006-01-26; 修订日期: 2006-06-01.

作者简介: 林豫柏(1980-), 男, 福建连城人, 系统工程专业硕士研究生.

*通讯联系人: 罗键(1954-), 男, 福建连城人, 教授, 博士生导师. 主要从事现代集成制造系统、系统控制与优化等研究. Email: jianluo@xmu.edu.cn

模式识别、回归分析和函数拟合等问题中,并且有一套坚实的理论基础。

文中阐述了支持向量机理论基础、-支持向量回归机算法。利用支持向量机算法具有全局最优、学习过程迅速、结构简单、构建决策原则同时获得一系列支持向量及推广能力强等优点,构建利用钢材生产过程中的化学成分,预测钢材最终产品的力学性能模型。

1 -支持向量回归机

-支持向量回归机相对应的原始最优问题^[3]为

$$\min_{w, b, \xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^l (\xi_i + \xi_i^*) \quad (1)$$

$$s. t. ((w \cdot x_i) + b) - y_i + \xi_i, i = 1, \dots, l \quad (2)$$

$$y_i - ((w \cdot x_i) + b) - y_i + \xi_i^*, i = 1, \dots, l \quad (3)$$

$$\xi_i, \xi_i^* \geq 0, i = 1, \dots, l \quad (4)$$

式中: $w \in R^n, \xi \in R^l, b \in R; (\cdot)$ 表示向量有 * 号和无 * 号两种情况。

问题(1) ~ (4) 的对偶问题

$$\min_{\alpha, \alpha^*} \frac{1}{2} \sum_{i,j=1}^l (\alpha_i^* - \alpha_i) (\alpha_j^* - \alpha_j) (x_i \cdot x_j) + \sum_{i=1}^l (\alpha_i^* + \alpha_i) - \sum_{i=1}^l y_i (\alpha_i^* - \alpha_i) \quad (5)$$

$$s. t. (\alpha_i - \alpha_i^*) = 0 \quad (6)$$

$$0 \leq \alpha_i, \alpha_i^* \leq \frac{C}{l}, i = 1, \dots, l \quad (7)$$

对偶问题式(5) ~ (7) 对输入 $x_i (i = 1, \dots, l)$ 的依赖关系仅仅体现在内积 $(x_i, x_j) (i, j = 1, \dots, l)$ 上, 因而其解 α^* 仅仅依赖于内积。利用这一特点, 通过引进核函数, 可把线性回归方法推广到处理非线性回归问题。

对于非线性回归, 首先使用一个非线性映射 ϕ 把数据映射到一个高维特征空间, 然后在高维特征空间进行线性回归。由于在上面的优化过程中只考虑到高维特征空间中的内积运算, 因此用一个核函数 $K(x, y)$ 代替 $\langle \phi(x), \phi(y) \rangle$ 就可以实现非线性回归。于是, 非线性回归的优化方程为最小化函数

$$w(\alpha, \alpha^*) = \frac{1}{2} \sum_{i,j=1}^l (\alpha_i^* - \alpha_i) (\alpha_j^* - \alpha_j) K(x_i, x_j) + \sum_{i=1}^l (\alpha_i^* + \alpha_i) - \sum_{i=1}^l y_i (\alpha_i^* - \alpha_i) \quad (8)$$

约束条件为式(6), (7)。

综上所述, -支持向量回归机算法可归纳为

- 1) 已知训练集 $T = \{(x_i, y_i), i = 1, \dots, l\}$, 其中, $x_i \in X \subset R^n, y_i \in Y \subset R, i = 1, \dots, l$;
- 2) 选取适的正数 C 和 γ , 以及适当的核函数;

3) 构造并求解最优化问题式(8), (6), (7), 得到最优解

$$\bar{w} = (\bar{w}_1, \bar{w}_2, \dots, \bar{w}_l, \bar{w}_{l+1}) \quad (9)$$

4) 构造决策函数

$$f(x) = \sum_{i=1}^l (\bar{w}_i^* - \bar{w}_i) K(x_i, x) + \bar{b} \quad (10)$$

其中, \bar{b} 按下列方式计算: 选择位于开区间 $(0, \frac{C}{l})$ 中的 \bar{w}_j 或 \bar{w}_k^* 。若选到的是 \bar{w}_j , 则

$$\bar{b} = y_j - \sum_{i=1}^l (\bar{w}_i^* - \bar{w}_i) (x_i, x_j) + \quad (11)$$

若选到的是 \bar{w}_k^* , 则

$$\bar{b} = y_k - \sum_{i=1}^l (\bar{w}_i^* - \bar{w}_i) (x_i, x_k) - \quad (12)$$

2 预测模型确定

2.1 核函数选择

选择适当的核函数关系到预测结果的好坏。文中选用的核函数是高斯径向基核函数 (Radial Basis Function, RBF)^[4]:

$$K(x_i, x_j) = \exp\left(-\gamma \|x_i - x_j\|^2\right), \gamma > 0 \quad (13)$$

2.2 最优参数 C, γ 的实现

在生成预测模型之前必须求出参数 C, γ 和 C , 参数 C, γ 和 C 对预测准确度有很大影响。常用的方法是把已知的训练集分成两部分: 一部分作为训练集, 其余部分作为测试集。根据训练集, 用被评价的分类算法求出决策函数, 用测试集测试所得的决策函数的准确率^[5-6]。交叉确认是该方法的改进。

k -折交叉确认把 l 个样本点随机的分成 k 个互补相交的子集, 即 k -折 S_1, \dots, S_k , 每折的大小相等。共进行 k 次训练与测试, 即对 $i = 1, 2, \dots, k$ 进行 k 次迭代, 第 i 次迭代的做法是, 选择 S_i 为测试集, 其余 $k-1$ 个集合的并为训练集, 算法根据训练集求出决策函数后, 即可对测试集 S_i 进行测试。

为获得最优参数 C, γ 和 C , 文中采用基于 5-折交叉确认的栅格搜索。将每一组 (C, γ) 进行 5-折交叉确认, 选出交叉确认误差最小的一组。该组合即为想要的最优参数组。实践表明, 使用指数变化序列的 C, γ 是非常实用的, 文中 C, γ 的指数变化序列为 $C = 2^{-4}, 2^{-3}, \dots, 2^4, \gamma = 2^{-8}, \dots, 1, \gamma = 2^{-8}, \dots, 2^{-1}$ 。 C, γ 分别取在变化序列不同值进行 5-折交叉确认, 比较交叉确认误差, 得到交叉确认误差最小的 C, γ 就是最优参数。

为了评价回归学习质量, 计算了均方误差 (MSE)^[7]。设延伸率的变换幅度

$$\nabla y = \max(y) - \min(y), \text{均方误差取}$$

$$MSE = \frac{1}{l} \sum_{i=1}^l \left(\frac{(y_i - f(x_i))}{\nabla y} \right)^2 \tag{14}$$

2.3 模型生成和预测

根据已获得的最优参数，和 C 值,生成预测模型

3 实验结果

实验训练样本采用三明钢铁集团公司生产的 HRB 335 产品的数据 , 总共 283 条记录 , 预测数据共 50 条记录 , 其部分数据见表 1 , 2.

表 1 HRB 335 延伸率部分样本数据

Tab.1 Sample data of HRB 335 extend rate

延伸率/ %	C	Mn	S	P	Si
30	0.20	1.39	0.039	0.028	0.50
29	0.22	1.43	0.027	0.034	0.50
28	0.20	1.40	0.030	0.023	0.46
...

表 2 HRB 335 延伸率部分预测数据

Tab.2 Forecast data of HRB 335 extend rate

延伸率/ %	C	Mn	S	P	Si
0	0.21	1.35	0.029	0.024	0.53
0	0.22	1.39	0.031	0.038	0.50
0	0.23	1.45	0.024	0.034	0.55
...

样本数据得出最优参数组合为 $\gamma = 0.125$, $\epsilon = 0.125$, $C = 4$, $MSE = 1.88$. 最优参数组合下预测值与实际值见图 1.

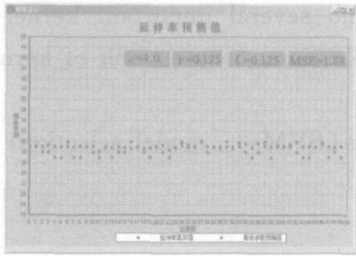


图 1 最优参数组合下预测值与实际值

Fig.1 Forecast value and practical value of optimal parameters combination

参考文献:

[1] Vladimir N Vapnik. An overview of statistical learning theory[J]. IEEE Transctions on Neural Networks , 1999 ,10 (5) : 988-999.

[2] Grace Wahba. Margin-like quantities and generalized approximate cross validation for support vector validation for support vector machines[J]. Neural Networks for Signal Processing , 1999 ,12 :20-23.

[3] 陶小龙. 基于支持向量机的股市预测[D]. 北京 :北京工业大学 ,2005.

[4] 邓乃杨. 数据挖掘的新方法 ——支持向量机[M]. 北京 :科学出版社 ,2004.

[5] 王靖华 ,何迪. 基于数据包字节长度的线性自回归和支持向量分类机的网络流量预测建模与分析[J]. 电脑应用 ,2005 (11) :133-135.

[6] 蒋琦庄 ,毅谢东. 基于 SVM 分类器的 SYN Flood 攻击检测规则生成方法的研究[J]. 计算机应用与软件 ,2005 (10) :83-86.

[7] 徐启华 ,杨端. 一种新的软间隔支持向量机分类算法[J]. 计算机工程与设计 ,2005 (9) :49-50.

4 与 Excel 回归方程预测比较

Excel 自带回归方程预测法是将样本导出 Excel ,选中样本后调用 Excel 函数得出回归方程 ,然后用回归方程进行预测. 为了与支持向量回归机作比较 ,样本采用三明钢铁集团公司的 HRB 335 产品的同样数据 ,共 283 条记录 ,部分样本数据见表 1. 所得回归方程为

延伸率 = 5.47 C + 30.09 Mn - 43.72 S -
2.57 P - 4.63 Si - 10.43 \tag{16}

基于支持向量机预测值、Excel 回归预测值与真实值见图 2.

通过比较表明 ,支持向量机预测比 Excel 回归预测具有更高的精度.

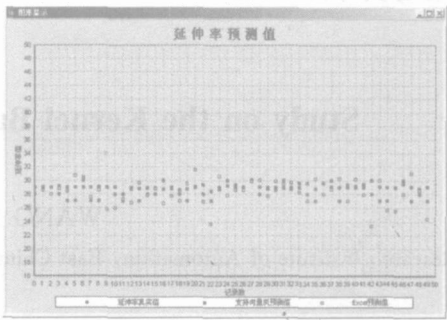


图 2 支持向量机预测值、Excel 回归预测值与真实值

Fig.2 Value of SVM, value of excel recursive forecast and real value

5 结 语

文中用基于支持向量回归的方法进行钢材力学性能预测 ,预测值与真实值相对误差小于 $\pm 8\%$,具有较高的预测精度. 该方法已在福建三明钢铁集团公司生产过程中得到实际应用 ,及时有效地进行了产品生产过程中的质量控制 ,从而提高了产品的质量.

(责任编辑 :杨 勇)